

JOURNAL OF MULTIVARIATE ANALYSIS 29, 127–136 (1989)

Diffusions of Perturbed Principal Component Analysis

MARIE-FRANCE BRU

*Université Paris Nord, 93 430 Villetaneuse, France**Communicated by the Editors*

We propose a stochastic differential equation approach to principal component analysis. We give the equations governing the spectrum of the square $B^T B$ of a $n \times p$ matrix of independent Brownian motions. We apply this result to P.C.A. of perturbed continuous data. © 1989 Academic Press, Inc.

A natural and usual way of studying stability in principal component analysis is to add a gaussian perturbation to data. But the resulting distributions are then very complicated [3, p. 193]. We have considered the problem a bit differently here by adding Brownian perturbations to the data. The first interest of the procedure is the possibility to make use of stochastic calculus methods. The second is to permit easy simulations.

Let $B(t)$ be a $n \times p$ matrix of independent Brownian motions, beginning in position $B(0) = D$, where $D = (d_{ij})$ is a $n \times p$ determinist matrix.

We wish to study the evolution of the eigenvalues and eigenvectors of the square matrix $B^T B$. We shall follow the method proposed by D. Williams in his note on Brownian motions of symmetric matrices [8].

1. RESULTS

THEOREM 1. *Let*

$$X(t) = B(t)^T B(t), \quad \text{where } B(t) = (b_{ij}(t)) \quad (1.1)$$

is a $n \times p$ -dimensional Brownian motion beginning in position $D = (d_{ij})$. We suppose D such that $X(0) = D^T D$ has p distinct eigenvalues:

$$\lambda_1(0) > \lambda_2(0) > \dots > \lambda_p(0).$$

Received October 21, 1987; revised February 22, 1988.

AMS 1980 subject classifications: 62H10, 60H10.

Key words and phrases: stochastic differential equations, sample covariance matrix, eigenvalues, principal component analysis.

Then at each time t , $X(t)$ will have, with probability 1, p distinct eigenvalues:

$$\lambda_1(t) > \lambda_2(t) > \dots > \lambda_p(t)$$

and the process $(\lambda_1, \lambda_2, \dots, \lambda_p)$ satisfies the stochastic differential equation

$$d\lambda_i(t) = 2 \sqrt{\lambda_i(t)} (dv_i(t)) + n dt + \sum_{k \neq i} \frac{\lambda_i(t) + \lambda_k(t)}{\lambda_i(t) - \lambda_k(t)} dt,$$

where v_1, v_2, \dots, v_p are independent Brownian motions.

Remark. It is easy to show that the elements $x_{ii}(t)$ of the diagonal of $X(t)$ are governed by the stochastic differential equations

$$dx_{ii}(t) = 2 \sqrt{x_{ii}(t)} (d\beta_i(t)) + n dt,$$

where $\beta_1, \beta_2, \dots, \beta_p$ are independent Brownian motions.

So as in the case treated by Dyson, Mc Kean [4], and Williams [8], the eigenvalues $\lambda_1, \lambda_2, \dots, \lambda_p$ behave here as the diagonal of X but subject to repulsion forces.

THEOREM 2. *With the same hypothesis as in Theorem 1, $X(t)$ can be diagonalized via an orthogonal transformation $H(t)$:*

$$H(t)^T X(t) H(t) = A(t) = \text{diag}(\lambda_i(t)) \quad (1.2)$$

in such a way that the $H(t)$ process is a continuous semimartingale in t . We have the equation

$$dH = H(dA + \frac{1}{2}(dA)(dA)),$$

where $A(t) = (a_{ij}(t))$ is a skew-symmetric $p \times p$ matrix, such that, for $i < j$,

$$da_{ij}(t) = \sqrt{\frac{\lambda_i(t) + \lambda_j(t)}{(\lambda_i(t) - \lambda_j(t))^2}} (d\beta_{ij}(t)).$$

Hence the eigenvectors satisfy the equation

$$\begin{aligned} dh_{ij}(t) = & \sum_{k \neq j} h_{ik}(t) \sqrt{\frac{\lambda_k(t) + \lambda_j(t)}{(\lambda_k(t) - \lambda_j(t))^2}} (d\beta_{kj}(t)) \\ & - \frac{1}{2} h_{ij}(t) \sum_{k \neq j} \frac{\lambda_k(t) + \lambda_j(t)}{(\lambda_k(t) - \lambda_j(t))^2} dt, \end{aligned}$$

where $\{\beta_{ij}, 1 \leq i < j \leq p\}$ is a family of Brownian motions all independent of one another and independent of the process $(\lambda_1, \lambda_2, \dots, \lambda_p)$.

2. SOME NOTATION

Let us write: d for the Ito differential.

Recall: if X and Y are matrix-valued semimartingales, the Ito rule for differentiating the product $X^T Y$ states:

$$d(X^T Y) = X^T(dY) + (dX)^T Y + (dX)^T (dY). \quad (2.1)$$

Remark. If we note ∂ for Stratonovich differential, we have

$$Y^T(\partial X) = Y^T(dX) + \frac{1}{2}(dY)^T (dX) \quad (2.2)$$

$$(\partial Y)^T X = (dY)^T X + \frac{1}{2}(dY)^T (dX);$$

(2.1) can then also be written

$$d(X^T Y) = X^T(\partial Y) + (\partial X)^T Y. \quad (2.3)$$

Notation. As in [8] if X and Y just differ by a finite-variation process, we shall write

$$dY \doteq dX.$$

3. PROOFS OF THEOREMS

(a) Let $\{b_{ij}, 1 \leq i \leq n, 1 \leq j \leq p\}$ be a family of independent Brownian motions. We have

$$(db_{ki})(db_{lj}) = \delta_{kl} \delta_{ij} dt.$$

The Ito rule (2.1) applied to (1.1) gives

$$dX = B^T(dB) + (dB)^T B + nI dt, \quad (3.1)$$

where I is the unit matrix of \mathbb{R}^p . Remark that

$$(B^T(dB))_{ij}((dB)^T B)_{kl} = x_{il} \delta_{jk} dt$$

gives

$$(dx_{ij})(dx_{kl}) = (x_{ik} \delta_{jl} + x_{il} \delta_{jk} + x_{jk} \delta_{il} + x_{jl} \delta_{ik}) dt,$$

so if $i \neq j$,

$$(dx_{ij})^2 = (x_{ii} + x_{jj}) dt$$

and if $i = j$,

$$(dx_{ii})^2 = 4x_{ii} dt.$$

The nondiagonal terms of $X(t)$ are martingales, and the diagonal terms are orthogonal semimartingales,

$$(dx_{ii})(dx_{jj}) = 4\sqrt{x_{ii}}\sqrt{x_{jj}}\delta_{ij}dt,$$

such that

$$x_{ii}(t) = x_{ii}(0) + nt + 2 \int_0^t \sqrt{x_{ii}(s)} (dv_i(s)),$$

where $\{v_i, 1 \leq i \leq p\}$ is a family of independent Brownian motions.

In particular, this equation shows that the x_{ii} terms of the diagonal are squares of Bessel processes (cf. [2, p. 224]).

(b) Let $\tau = \inf\{s/\lambda_i(s) = \lambda_j(s) \text{ for a couple } (i, j)\}$; for $t < \tau$, the eigenvalues $(\lambda_i(t))_{1 \leq i \leq p}$ of semimartingale matrices $X(t)$ being C^∞ functions of $\{x_{ij}(t)\}$, are semimartingales (cf. [2, p. 66]) and we can choose a family $H(t)$ of orthogonal matrices, continuous in t , which diagonalises $X(t)$:

$$H(t)^T X(t) H(t) = \Lambda(t) = (\text{diag } \lambda_i(t)). \quad (1.2)$$

The coefficients of $H(t)$ being C^∞ functions of (x_{ij}) and (λ_i) , are also semimartingales on $t < \tau$ (see also D. Williams [8]).

Let us now apply (2.1) to (1.2),

$$\begin{aligned} d\Lambda &= (dH)^T X H + H^T (dX) H + H^T X (dH) + (dH)^T (dX) H \\ &\quad + (dH)^T X (dH) + H^T (dX) (dH) \end{aligned} \quad (3.2)$$

and introduce some auxiliary matrices:

i. $A(t)$ ("Logarithm" of H) such that

$$A(0) = 0, \quad dA = H^T (dH) + \frac{1}{2}(dH)^T (dH) \quad (\text{or } dA = H^{-1}(\partial H)),$$

at each time t , $A(t)$ is skew-symmetric, and

$$dH = H((dA) + \frac{1}{2}(dA)(dA)) \quad (\text{or } dH = H(\partial A)). \quad (3.3)$$

ii. $d\Gamma = \frac{1}{2}(dA)(dA)$, $\Gamma(0) = 0$ is a finite -variation process, and (3.3) gives

$$dH = H(dA + d\Gamma). \quad (3.4)$$

iii. $d\Phi = H^T(dX)H(dA)$; $d\mu = (dA)^T \Lambda(dA)$. With all these new matrices, (3.2) is

$$\begin{aligned} d\Lambda &= H^T(dX)H + ((dA)^T \Lambda + \Lambda(dA)) \\ &\quad + ((d\Gamma)^T \Lambda + \Lambda(d\Gamma)) + d\Phi + d\Phi^T + d\mu. \end{aligned} \quad (3.5)$$

(c) *Proof of Theorem 1.* i. The diagonal terms of (3.5) give for the martingale part,

$$d\lambda_i \doteq \sum_{k,l} h_{ki} h_{li} (dx_{kl}) \doteq 2 \sum_{k,l,r} h_{ki} h_{li} b_{rk} (db_{rl}), \quad (3.6)$$

and for the finite-variation part,

$$\begin{aligned} &[\text{Finite-variation part of } d\lambda_i] \\ &= [\text{Finite-variation part of } (H^T(dX)H)_{ii}] \\ &\quad + 2\lambda_i (d\gamma_{ii}) + 2d\phi_{ii} + d\mu_{ii}. \end{aligned} \quad (3.7)$$

ii. The nondiagonal terms of (3.5) give:

$$\begin{aligned} (\lambda_j - \lambda_i)(da_{ij}) &= \sum_{k,l} h_{ki} h_{lj} (dx_{kl}) \\ &\quad + (\lambda_i + \lambda_j)(d\gamma_{ij}) + d\phi_{ij} + d\phi_{ji} + d\mu_{ij} \end{aligned} \quad (3.8)$$

so that if $i \neq j$ and $m \neq k$,

$$\begin{aligned} &\frac{(\lambda_j - \lambda_i)(\lambda_m - \lambda_k)(da_{ij})(da_{km})}{dt} \\ &= \lambda_i(\delta_{jm} \delta_{ik} + \delta_{jk} \delta_{im}) + \lambda_j(\delta_{im} \delta_{jk} + \delta_{ik} \delta_{jm}). \end{aligned} \quad (3.9)$$

iii. If $t < \tau$,

$$\begin{aligned} d\Gamma &= \text{diag}(d\gamma_{ii}) \\ &= \text{diag}\left(\frac{1}{2} \sum_k (da_{ik})(da_{ki})\right) \\ &= \text{diag}\left(-\frac{1}{2} \sum_{k \neq i} \frac{\lambda_i + \lambda_k}{(\lambda_i - \lambda_k)^2} dt\right) \end{aligned} \quad (3.10)$$

$$\begin{aligned} d\mu &= \text{diag}(d\mu_{ii}) \\ &= \text{diag}\left(\sum_k (da_{ki})^2 \lambda_k\right) \\ &= \text{diag}\left(\sum_{k \neq i} \frac{\lambda_i + \lambda_k}{(\lambda_i - \lambda_k)^2} \lambda_k dt\right), \end{aligned} \quad (3.11)$$

with the help of formula (3.8), we have

$$(\lambda_j - \lambda_q)(da_{qj})(dx_{ki}) = (\lambda_j + \lambda_q)(h_{kj}h_{iq} + h_{ij}h_{kq}) dt, \quad (3.12)$$

so

$$d\Phi = \text{diag}(d\phi_{ii}) = \text{diag} \left(\sum_{q \neq i} \frac{\lambda_i + \lambda_q}{\lambda_i - \lambda_q} dt \right). \quad (3.13)$$

It now follows from (3.1), (3.7), (3.10), (3.11), and (3.13), that the finite-variation part of $d\lambda_i$ is

$$[\text{Finite-variation part of } d\lambda_i] = \left(n + \sum_{k \neq i} \frac{\lambda_i + \lambda_k}{\lambda_i - \lambda_k} \right) dt. \quad (3.14)$$

If we now remark with (3.6) that

$$(d\lambda_i)(d\lambda_j) = 4 \sqrt{\lambda_i} \sqrt{\lambda_j} \delta_{ij} dt \quad (3.15)$$

and if we suppose the time of the first collision is infinite ($\tau = +\infty$ a.s.), Theorem 1 is then proved.

(d) *Proof of Theorem 2.* For $t < \tau$, matrices $d\Gamma$, $d\Phi$, $d\mu$, [Finite-variation part of $H^T(dX)H$] are all diagonal, and the (3.8) relation shows that dA is a matrix-valued martingale,

$$da_{ij} = \frac{1}{\lambda_j - \lambda_i} \sum_{k,l,r} h_{ki} h_{lj} (b_{rk}(db_{rl}) + b_{rl}(db_{rk})) \quad (3.16)$$

with

$$(da_{ij})^2 = \frac{\lambda_i + \lambda_j}{(\lambda_i - \lambda_j)^2} dt, \quad da_{ij} = \sqrt{\frac{\lambda_i + \lambda_j}{(\lambda_i - \lambda_j)^2}} (d\beta_{ij}), \quad (3.17)$$

and the eigenvectors which are given by

$$dH = (dh_{ij}) = H(dA + \frac{1}{2}(dA)(dA))$$

are

$$dh_{ij} = \sum_{k \neq j} h_{ik} \sqrt{\frac{\lambda_j + \lambda_k}{(\lambda_j - \lambda_k)^2}} (d\beta_{kj}) - \frac{1}{2} h_{ij} \sum_{k \neq j} \frac{\lambda_k + \lambda_j}{(\lambda_k - \lambda_j)^2} dt, \quad (3.18)$$

where $\{\beta_{ij}, 1 \leq i < j \leq p\}$ is a family of independent Brownian motions.

We also have from (3.6) and (3.8), for all i, k, m ,

$$0 = (d\lambda_i)(da_{km}) = 2 \sqrt{\lambda_i} (dv_i) \sqrt{\frac{\lambda_k + \lambda_m}{(\lambda_k - \lambda_m)^2}} (d\beta_{km})$$

and, as the σ -algebra, $\sigma\{\lambda_i, 1 \leq i \leq p\}$ is included in the σ -algebra $\sigma\{v_i, 1 \leq i \leq p\}$. Theorem 2 is proved.

(e) *The impossibility of collision.* The demonstrations of Theorems 1 and 2 would not be complete if we did not prove that the time of the first collision between two eigenvalues is infinite ($\tau = +\infty$ a.s.).

Here again we shall follow the method proposed by Williams [8]: Find a real function U , defined on $D = \{(x_1, x_2, \dots, x_p) \in \mathbb{R}^p / x_1 > x_2 > \dots > x_p\}$ with continuous derivatives and such that $\mathcal{U}(t) = U(\lambda_1(t), \lambda_2(t), \dots, \lambda_p(t))$ is a local martingale, infinite when two eigenvalues collide. We have seen that

$$d\lambda_i = dm_i + \psi_i dt,$$

where m_i is a martingale.

Apply Ito's formula to $\mathcal{U}(t)$, to obtain

$$d\mathcal{U} = \sum_i \frac{\partial U}{\partial \lambda_i} (dm_i) + \sum_i \frac{\partial U}{\partial \lambda_i} \psi_i dt + \frac{1}{2} \sum_{i,j} \frac{\partial^2 U}{\partial \lambda_i \partial \lambda_j} (d\lambda_i)(d\lambda_j) \quad (3.19)$$

as $(d\lambda_i)(d\lambda_j) = 4\lambda_i \delta_{ij} dt$, $\mathcal{U}(t)$ is a local martingale if the finite-variation part of (3.19) is null:

$$\sum_i \frac{\partial U}{\partial \lambda_i} \left(n + \sum_{k \neq i} \frac{\lambda_i + \lambda_k}{\lambda_i - \lambda_k} \right) + 2 \sum_i \frac{\partial^2 U}{\partial \lambda_i^2} \lambda_i = 0 \quad (3.20)$$

For $p=2$, (3.20) is

$$\frac{\partial U}{\partial x} \left(n + \frac{x+y}{x-y} \right) + \frac{\partial U}{\partial y} \left(n + \frac{y+x}{y-x} \right) + 2 \left(\frac{\partial^2 U}{\partial x^2} x + \frac{\partial^2 U}{\partial y^2} y \right) = 0 \quad (3.21)$$

on $D = \{(x, y) \in \mathbb{R}^2 / x > y\}$.

It is natural to look for solutions $U(x, y)$ of (3.21) which are functions of the difference $x - y$, null if there is a collision.

So let

$$U(x, y) = f(x - y) \quad \text{and} \quad z = x - y.$$

(3.21) becomes

$$zf''(z) + f'(z) = 0$$

whose solution on $]0, \infty[$ is

$$f(z) = a \operatorname{Log} z + b,$$

so

$$U(x, y) = \text{Log}(x - y)$$

is a solution of (3.21). It is then natural and easy to show that

$$U(x_1, x_2, \dots, x_p) = \sum_{i < j} \text{Log}(x_i - x_j)$$

is a solution of (3.20) on D . Consequently,

$$\mathfrak{U}(t) = U(\lambda_1(t), \lambda_2(t), \dots, \lambda_p(t)) = \sum_{i < j} \text{Log}(\lambda_i(t) - \lambda_j(t))$$

which is continuous on $[0, \tau[$ and such that

$$\lim_{t \uparrow \tau} \mathfrak{U}(t) = -\infty$$

is also a local martingale on $[0, \tau[$, so it is a time-transformation of Brownian motion.

We now can follow the argument of Mc Kean [4, p. 47], and Williams [8]. If $\zeta(t)$ is on $[0, \tau[$, the inverse function of

$$\langle U \rangle(t) = \int_0^t \sum_i \sum_{j \neq i} \left(\frac{1}{\lambda_i(s) - \lambda_j(s)} \right)^2 ds,$$

$B(t) = U(\zeta(t))$ is a Brownian motion on $[0, \langle U \rangle(\tau)[$ [2, p. 92]; $\tau < +\infty$ would imply

$$\lim_{t \uparrow \langle U \rangle(\tau)} B(t) = \lim_{t \uparrow \tau} \mathfrak{U}(t) = -\infty$$

which is impossible if $B(t)$ is a Brownian motion. Thus $\tau = +\infty$ a.s., and the result follows.

4. APPLICATION TO SPECTRAL ANALYSIS OF BROWNIAN SAMPLE VARIANCE-COVARIANCE MATRICES

Let

$$\tilde{B}(t) = (b_{ij}(t) - \bar{b}_j(t)), \quad \text{where} \quad \bar{b}_j(t) = \sum_q b_{jq}(t);$$

we shall now apply the preceding results to spectral analysis of the matrices

$$S(t) = \tilde{B}(t)^T \tilde{B}(t).$$

Such matrices are classical in principal component analysis.

This case can easily be deduced from the precedent by an orthogonal transformation. Let $\Theta = (\theta_{ij})$ be an orthogonal matrix of \mathbb{R}^n , whose last line is $(1/\sqrt{n})(\theta_{ni} = 1/\sqrt{n}, 1 \leq i \leq p)$. Such a matrix transforms B in $B' = \Theta B$, with the last line $b'_{nj} = \sqrt{n} \bar{b}_j$. As Θ is constant

$$dB' = \Theta(dB) \quad \text{and} \quad (dB')^T (dB) = nI dt,$$

so B' is a matrix of independent Brownian motions such that $B'(0) = \Theta D = D'$.

Let us now note \bar{B}' , the matrix B' from which the last line has been drawn out; we then have

$$S(t) = \bar{B}'(t)^T \bar{B}'(t) = \left(\sum_k (b_{ki}(t) - \bar{b}_i(t))(b_{kj}(t) - \bar{b}_j(t)) \right) = \bar{B}(t)^T \bar{B}(t).$$

If we now apply Theorems 1 and 2 to $S(t)$, we obtain the same results, with just one difference in the stochastic differential equations which govern the eigenvalues, n becomes $n-1$:

$$d\lambda_i(t) = \sqrt{\lambda_i(t)} (dv_i(t)) + (n-1) dt + \sum_{k \neq i} \frac{\lambda_i(t) + \lambda_k(t)}{\lambda_i(t) - \lambda_k(t)} dt.$$

Here again the eigenvalues $\lambda_1, \lambda_2, \dots, \lambda_p$ behave as the diagonal of S , but subject to repulsion forces.

5. DISTRIBUTIONS

Let us now give some remarks on the distributions of the X and S semimartingales.

For a fixed t , $\{b_{ij}(t), 1 \leq i \leq n, 1 \leq j \leq p\}$ are independent gaussian variables ($b_{ij}(t) \sim \eta(d_{ij}, t)$), so $X(t) = B(t)^T B(t)$ has a noncentral Wishart distribution with noncentrality parameter $D^T D$. The density of such a distribution can only be expressed by generalised Bessel functions. Yet, let us notice that the diagonal terms $x_{ii}(t)$ (which are governed by the equations: $dx_{ii} = 2\sqrt{x_{ii}}(db_i) + n dt$) have a noncentral $t\chi^2(n, x_{ii}(0)/t)$ distribution (cf. Rao [5, p. 165]), for the classical calculus by multivariate analysis, or [2, p. 224] for the same result obtained by stochastic calculus). The distribution of the eigenvalues of such a random matrix is known, but it is also very intricate [3, p. 193].

The same remarks hold for $S(t) = \bar{B}(t)^T \bar{B}(t)$ as this case results from the former by an orthogonal transformation. Maybe we could expect more information about those distributions from the stochastic differential equations we have obtained. Yet we can remark that by discretisation of

the above differential equations, it is possible to simulate the distributions of the eigenvalues and eigenvectors, without any additional diagonalisation on perturbed data (cf. [1]).

6. SPECTRAL ANALYSIS OF SAMPLE CORRELATION MATRICES

In principal component analysis we often have to study the spectrum of sample correlation matrices, such as,

$$R(t) = \hat{B}(t)^T \hat{B}(t), \quad R(0) = \hat{D}^T \hat{D},$$

where

$$\hat{B}(t) = \left(\frac{b_{ij}(t) - \bar{b}_j(t)}{\sqrt{s_{ij}(t)}} \right) = (\rho_{ij}(t)), \quad \hat{D} = \left(\frac{d_{ij} - \bar{d}_j}{\sqrt{s_{ij}(0)}} \right) = (\rho_{ij}(0)).$$

The same method of calculus gives, under noncollision hypothesis, the stochastic differential equations which govern the motion of the eigenvalues and eigenvectors of $R(t)$. Those equations are more complicated and are given in [1, pp. 207, 212].

7. REMARK

The paper [5] treats the case of Dynkin's Brownian motions of ellipsoids with analogous methods. In particular, the authors derive the corresponding stochastic differential equations and obtain a noncollision theorem of eigenvalues (see also [7, IV.36]).

REFERENCES

- [1] BRU, M. F. (1987). Thèse de 3ème cycle, Université Paris Nord.
- [2] IKEDA, N. AND WATANABE, S. (1981). *Stochastic Differential Equations and Diffusion Processes*. North-Holland, Amsterdam.
- [3] JOHNSON, N. L., AND KOTZ, S. (1972). *Distributions in Statistics: Continuous Multivariate Distributions*. Wiley, New York.
- [4] MCKEAN, H. P. (1969). *Stochastic Integrals*. Academic Press, New York.
- [5] NORRIS, J. K., ROGERS, L. C. G., AND WILLIAMS, D. (1986). Brownian motions of ellipsoids. *Trans. Amer. Math. Soc.* **294** 757–765.
- [6] RAO, C. R. (1965). *Linear Statistical Inference and Its Applications*. Wiley, New York.
- [7] ROGERS, L. C. G., AND WILLIAMS, D. (1987). *Diffusions, Markov Processes and Martingales*, Vol. 2. *Ito Calculus*. Wiley, New York.
- [8] WILLIAMS, D. (1985). A note on Brownian motions and symmetric matrices. Preprint, University of Cambridge.